# K-MEANS CLUSTERING

Let $x_1, ..., x_n \in \mathbb{R}^d$.

Recall that the goal of clustering is to
assign th data to disjoint subsets called

(A) _____ so that points in the same

cluster are more similar than points in

different clusters.

Therefore, at the heart of every clustering
algorithm is a notion of _____.
Often it is more convenient to work with

a _____.

# Dissimilarity

A dissimilarity matrix is an $n \times n$ matrix

$$D = \left[ d_{ij} \right]_{i,j=1}^{n}$$

which has the following properties
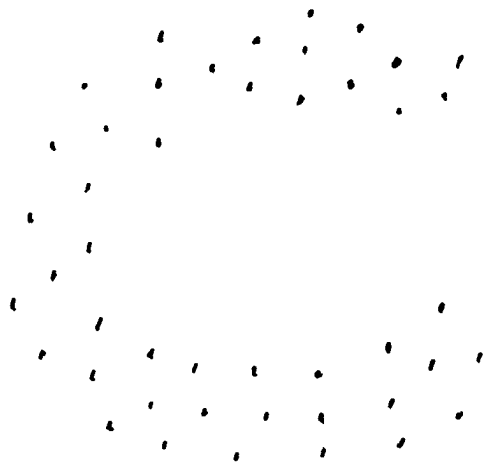
- $d_{ii} = 0$
- $d_{ij} = d_{ji}$
- $d_{ij} \geq 0$

Conceptually, if $x_i$ is more similar to $x_j$ than to $x_k$, then

Ⓑ

A dissimilarity matrix need not satisfy the triangle inequality:

# Examples

Ⓒ

- Euclidean distance

- Squared Euclidean distance

- kNN - based distance

# K-means criterion

A __cluster map__ is a function

$$C: \{1, 2, \ldots, n\} \longrightarrow \{1, 2, \ldots, K\}$$

that partitions the data into $K$ clusters.

In K-means clustering we

- assume $K$ is known (more on this later)

- adopt the squared Euclidean distance as a dissimilarity

(D)
$$d_{ij} =$$

- seek to minimize the _____ _____

_____

$$W(c) =$$

where

$$n_k =$$

# Algorithm

(E) The K-means criterion is a _____ optimization problem. The number of possible cluster maps $C$ is

$$\frac{1}{K!} \sum_{k=1}^{K} (-1)^{K-k} \binom{K}{k} k^n \qquad \text{(Jain \& Dubes, 1988)}$$

$$\begin{cases} = 34,105 & \text{if } n=10, \; K=4 \\ \\ \approx 10^{10} & \text{if } n=19, \; K=4 \end{cases}$$

There is no known efficient search strategy for this space. Therefore we resort to an iterative, suboptimal algorithm.

## Exercise | Show that

$$W(C) = \sum_{k=1}^{K} \sum_{i: C(i)=k} \| x_i - \bar{x}_k \|^2$$

where

$$\bar{x}_k := \frac{1}{n_k} \sum_{i: C(i)=k} x_i$$

## Solution

$$W(C) = \frac{1}{2} \sum_{k=1}^{K} \frac{1}{n_k} \sum_{i:C(i)=k} \sum_{j:C(j)=k} \underbrace{\left\| x_i - \bar{x}_k - (x_j - \bar{x}_k) \right\|^2}$$
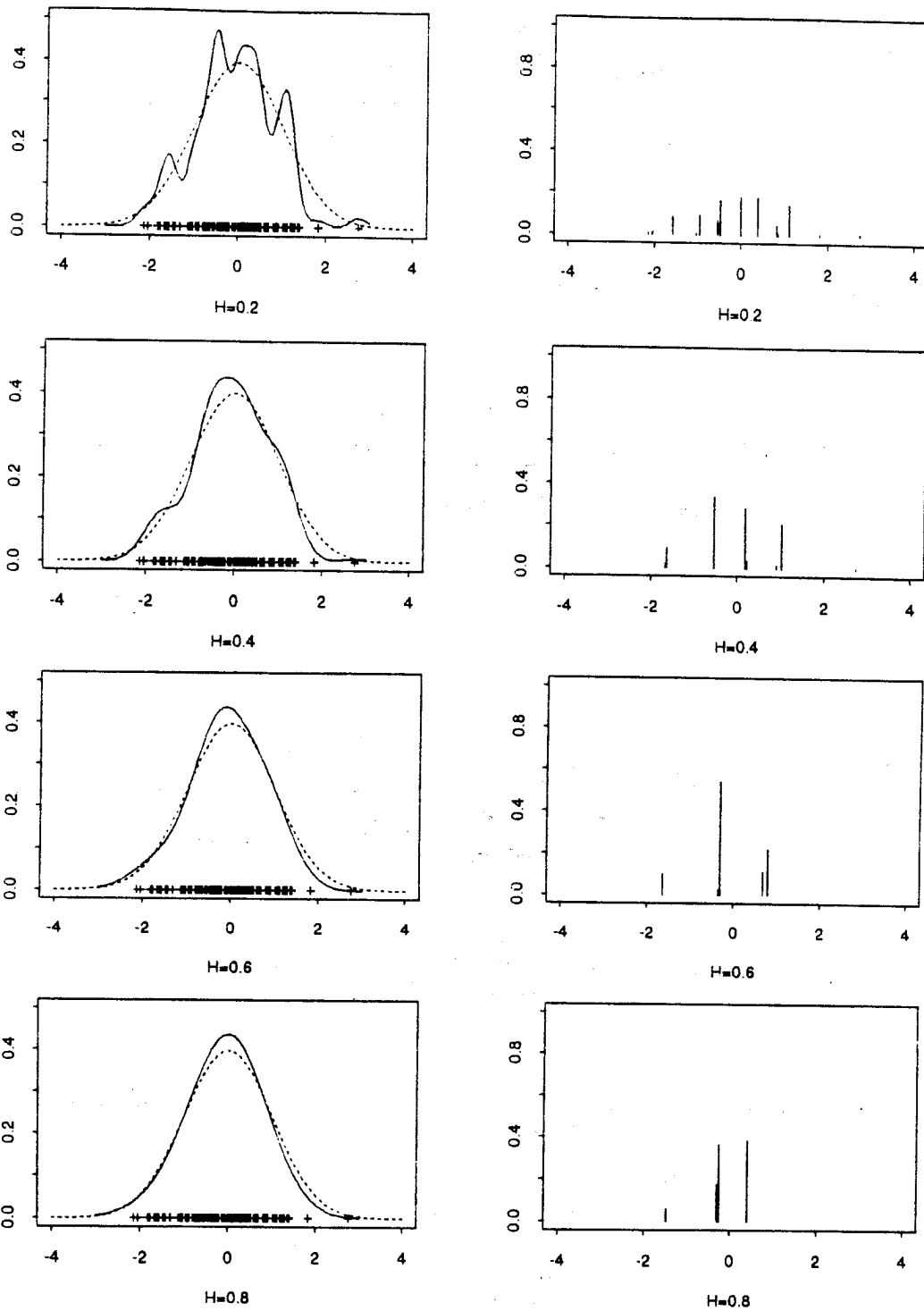
$$\langle x_i - \bar{x}_k - (x_j - \bar{x}_k), \; x_i - \bar{x}_k - (x_j - \bar{x}_k) \rangle$$

$$= \| x_i - \bar{x}_k \|^2 - 2(x_i - \bar{x}_k)^T (x_j - \bar{x}_k) + \| x_j - \bar{x}_k \|^2$$

$$= \frac{1}{2} \sum_{k=1}^{K} \frac{1}{n_k} \left[ \sum_{i:C(i)=k} \sum_{j:C(j)=k} \| x_i - \bar{x}_k \|^2 \right.$$

$$- 2 \sum_{i:C(i)=k} \sum_{j:C(j)=k} (x_i - \bar{x}_k)^T (x_j - \bar{x}_k)$$

$$\left. + \sum_{i:C(i)=k} \sum_{j:C(j)=k} \| x_j - \bar{x}_k \|^2 \right]$$

$$= \frac{1}{2} \sum_{k=1}^{K} \frac{1}{n_k} \left[ n_k \cdot \sum_{i:C(i)=k} \| x_i - \bar{x}_k \|^2 \right.$$

$$\left. + n_k \sum_{j:C(j)=k} \| x_j - \bar{x}_k \|^2 \right]$$

$$= \sum_{k=1}^{K} \sum_{i:C(i)=k} \| x_i - \bar{x}_k \|^2$$

**Figure 2-1-1 (A)** : LSMDE for $N(0,1)$ with positive weights ($n$=100). The solid line is LSMDE and the dotted line is $N(0,1)$.

32

G) 1) $m_k^* =$

2) $C^*(i) =$

---

### K-means Clustering Algorithm

Initialize $\bar{x}_k$, $k = 1, \ldots, K$

Repeat

- $C(i) =$

- $\bar{x}_k =$

Until clusters don't change

H)

---

Remarks

- The algorithm is typically initialized by setting each $\bar{x}_k$ to be a _random_ data point

- Since the algorithm often finds a local min, several random initializations are recommended.
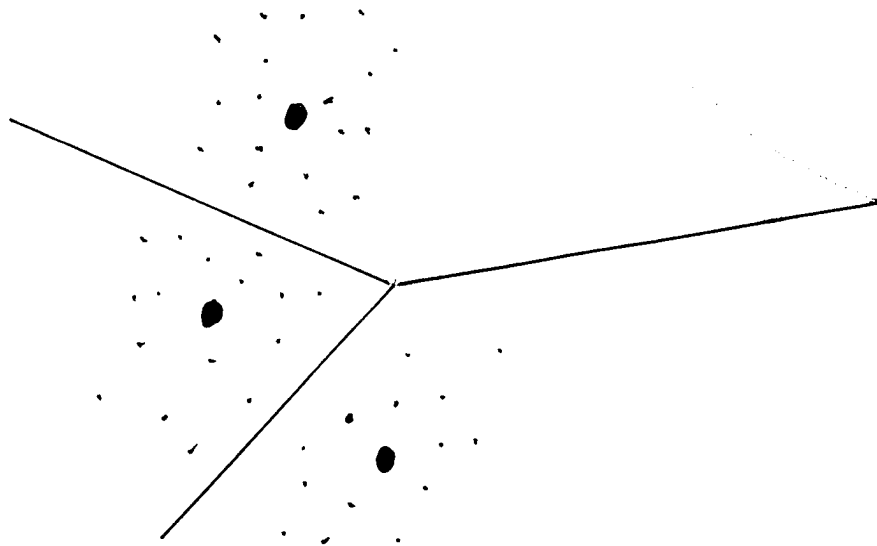
# Cluster Geometry

Clusters are "nearest neighbors"

① regions or _____ cells

defined with respect to the cluster means.

Therefore the cluster boundaries are

_____



K = 3
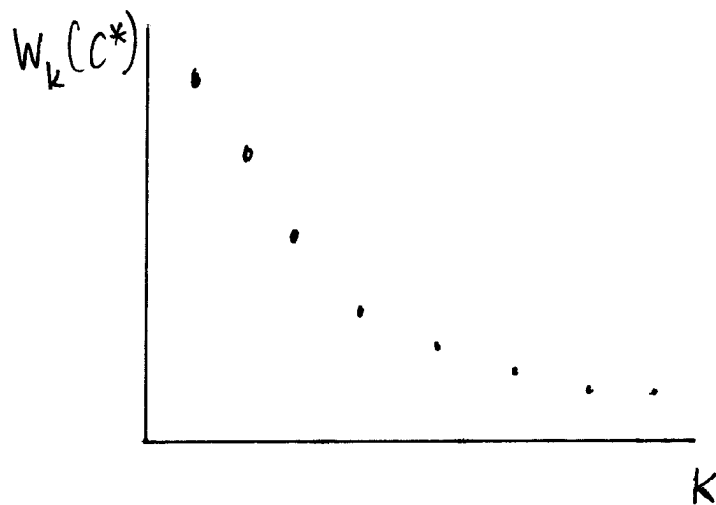
K-means will fail if clusters are

_____ .

## Model selection

How to choose $K$?

If $W_K(C^*)$ is the within-cluster scatter based on $K$ clusters, we have a plot like this



If the "right" number of clusters is $K^*$, we expect

- for $K < K^*$, $W_K(C^*) - W_{K-1}(C^*)$ will be <u>large</u>

- for $K > K^*$, $W_K(C^*) - W_{K-1}(C^*)$ will be <u>small</u>

This suggests choosing $K$ near the "knee" of the curve.

$\boxed{\text{Key}}$ A. clusters, similarity, dissimilarity

B. $\quad d_{ij} < d_{ik}$

$\quad d_{ij} + d_{jk} \not\geq d_{ik}$

C. $\bullet \ \|x - y\| = \left( \sum_{j=1}^{d} (x^{(j)} - y^{(j)})^2 \right)^{\frac{1}{2}}$

$\bullet \ \|x - y\|^2$

$\bullet$ length of shortest path on $k$-nearest neighbor graph, for some $k$

D. $\quad d_{ij} = \|x_i - x_j\|^2$, within cluster scatter

$$W(C) = \frac{1}{2} \sum_{k=1}^{K} \sum_{i:C(i)=k} \left[ \underbrace{\frac{1}{n_k} \sum_{j:C(j)=k} \|x_i - x_j\|^2}_{\text{avg. dissim. to points in same cluster}} \right]$$

$$n_k = \sum_{i=1}^{n} \mathbb{1}_{\{C(i)=k\}}$$

E. Combinatorial          F. $\overline{x}_k$

G. $\quad m_k^* = \dfrac{1}{n_k} \sum\limits_{i:C(i)=k} x_i$

$\quad C^*(i) = \underset{k}{\arg\min} \; \| x_i - \quad \|$

H. $\quad C(i) = \underset{k}{\arg\min} \; \| x_i - \bar{x}_k \|$

$\quad \bar{x}_k = \dfrac{1}{n_k} \sum\limits_{i:C(i)=k} x_i$

I. Voronoi , hyperplanes , nonconvex